



Πανεπιστήμιο Πειραιώς

Σχολή Τεχνολογιών Πληροφορικής και Τηλεπικοινωνιών

Τμήμα Ψηφιακών Συστημάτων

ΤΙΤΛΟΣ ΕΡΓΑΣΙΑΣ :

Απαλλακτική εργασία μαθήματος (Εαρ. Εξ. 2020-21)

ΣΤΟΙΧΕΙΑ ΦΟΙΤΗΤΗ :

- ❖ ΤΣΟΥΦΗΣ ΘΕΟΔΩΡΟΣ
- ❖ ΑΜ: Ε17155
- ❖ Εξάμηνο σπουδών: 8

Εργασία στο μάθημα «Διακυβέρνηση Πληροφοριακών Συστημάτων»

Επιβλέπουσα: κ.Μαυρογιώργου Αργυρώ.

ΠΙΝΑΚΑΣ ΠΕΡΙΕΧΟΜΕΝΩΝ

1 ΘΕΜΑΤΟΛΟΓΙΑ ΕΦΑΡΜΟΓΗΣ ΑΝΑΛΥΣΗΣ ΔΕΔΟΜΕΝΩΝ.	3
2.ΑΡΧΙΤΕΚΤΟΝΙΚΗ ΕΦΑΡΜΟΓΗΣ	4
3. ΕΓΧΕΙΡΙΔΙΟ ΧΡΗΣΗΣ ΕΦΑΡΜΟΓΗΣ	5
4. ΑΠΟΤΕΛΕΣΜΑΤΑ ΕΦΑΡΜΟΓΗΣ	10

ΕΝΟΤΗΤΑ 1: Θεματολογία Εφαρμογής Ανάλυσης Δεδομένων

Η εφαρμογή αυτή μαζεύει live streaming δεδομένα από το twitter . Τα δεδομένα που μαζεύει (όταν λέμε δεδομένα εννοούμε tweets χρηστών) περιέχουν κάποιες από τις ακόλουθες λέξεις κλειδιά → covid-19 , vaccine , vaccination , vaccinated , AstraZeneca , Pfizer .

Σκοπός της εφαρμογής είναι από τα live streaming δεδομένα που θα συλλέξει (τα οποία είναι 2000) , να κάνει μια ανάλυση και να μας δείξει πόσα από τα μηνύματα αυτά περιέχουν αρνητικές λέξεις-εκφράσεις . Δηλαδή πόσα από τα δεδομένα περιέχουν όχι και τόσο θετικά μηνύματα ή σκέψεις σχετικά με τον εμβολιασμό κατά του ιού Covid-19 . Η εφαρμογή αυτή θέλει να αποδείξει την αρνητικότητα-προβληματισμό-ανησυχία του κόσμου που επικρατεί έναντι του εμβολιασμού κατά του Covid-19 .

Η εφαρμογή αρχικά συλλέγει τα δεδομένα , μετά εμφανίζει κάποια στατιστικά στοιχεία και τέλος χωρίζει τα δεδομένα σε 2 ομάδες . Δεδομένα με θετικά μηνύματα ή σκέψεις σχετικά με τον εμβολιασμό κατά του ιού Covid-19 και δεδομένα με αρνητικά μηνύματα ή σκέψεις σχετικά με τον εμβολιασμό κατά του ιού Covid-19 .

ΕΝΟΤΗΤΑ 2: Αρχιτεκτονική Εφαρμογής

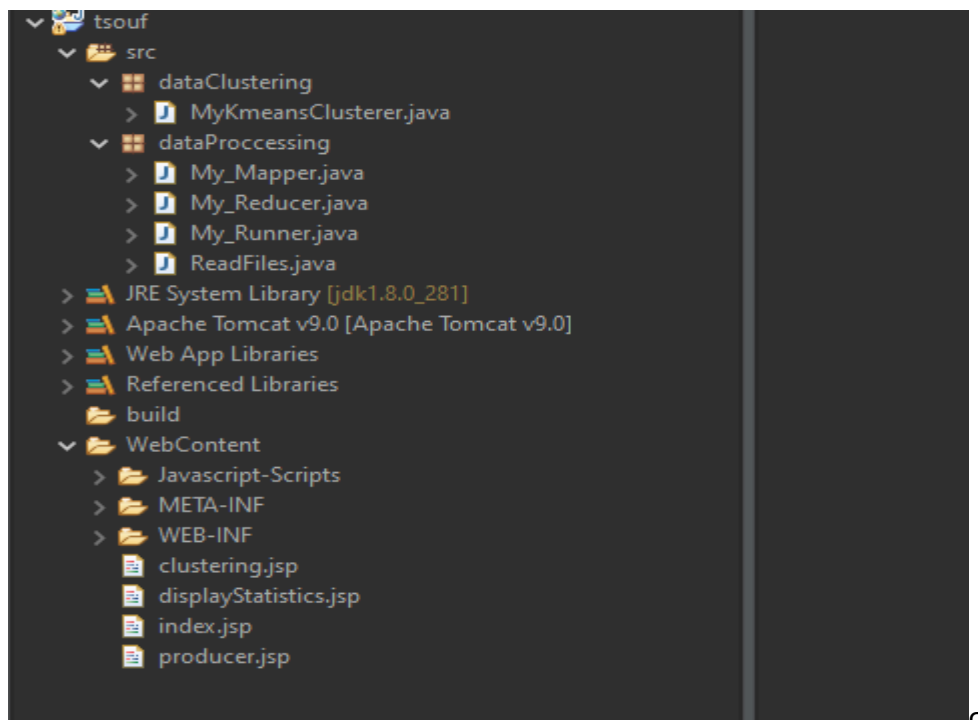
Η εφαρμογή έχει δομηθεί με jsp ,html , css ,bootStrap, javascript .Τα εργαλεία που χρησιμοποιήθηκαν για την συλλογή και επεξεργασία των δεδομένων είναι τα KAFKA , FLUME , HDFS, MapReduce , Mahout .

Πιο συγκεκριμένα έχουμε μία Java-Web Application εφαρμογή που όσον αφορά το front-end χρησιμοποιεί html , css ,bootStrap και javascript ενώ στο back-end χρησιμοποιείται η Java .

Η java στο backend χρησιμοποιεί τα εργαλεία KAFKA , FLUME , HDFS, MapReduce , Mahout έτσι ώστε να συλλέξει και να επεξεργαστεί τα live streaming δεδομένα .

Σε αρχικό στάδιο η εφαρμογή με τη βοήθεια του KAFKA και του FLUME συλλέγει τα δεδομένα και τα γράφει σε ένα φάκελο στο HDFS .Μετά εφαρμόζει ένα MapReduce job στα συλλεχθέντα δεδομένα με σκοπό να εμφανίσει κάποια στατιστικά .Μετά με την βοήθεια του Mahout η εφαρμογή χρησιμοποιεί τα συλλεχθέντα δεδομένα από το HDFS και τα προεπεξεργάζεται .Κατά την προεπεξεργασία των δεδομένων (tweets) διαχωρίζονται τα μηνύματα που έχουν αρνητική χροιά από τα υπόλοιπα. Τέλος εφαρμόζει στα προεπεξεργασμένα δεδομένα τον αλγόριθμο K-Means , με σκοπό να δούμε πόσα τελικά από τα 2000 tweets περιείχαν κάτι αρνητικό σχετικά με τον εμβολιασμό κατά του Covid-19 .

Ακολουθεί η δομή του Java-Web Application.

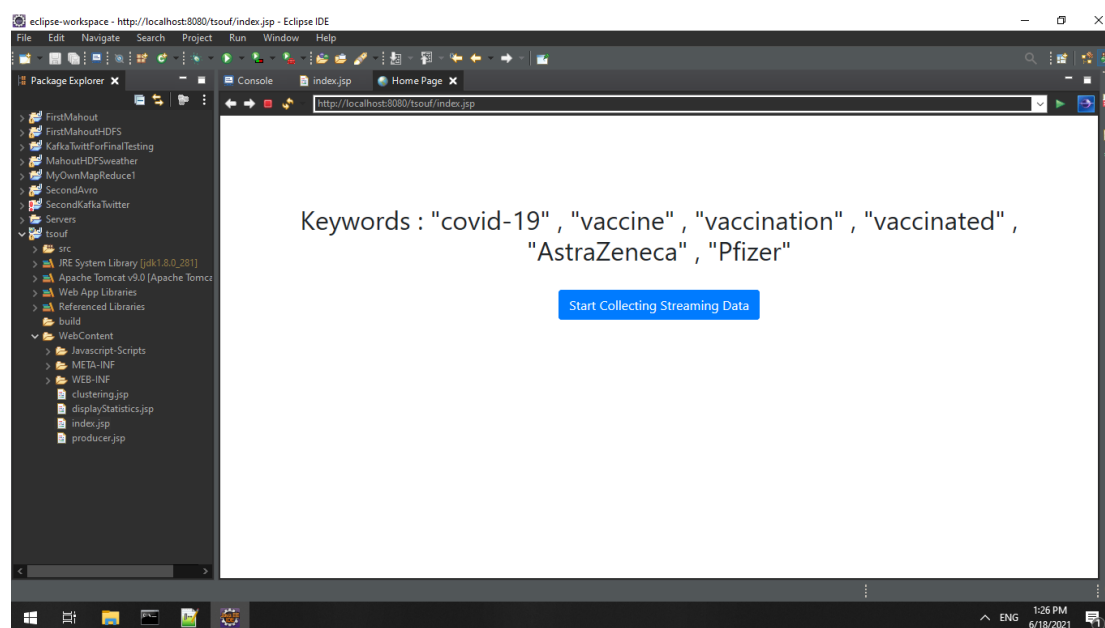


ΕΝΟΤΗΤΑ 3: Εγχειρίδιο Χρήσης Εφαρμογής

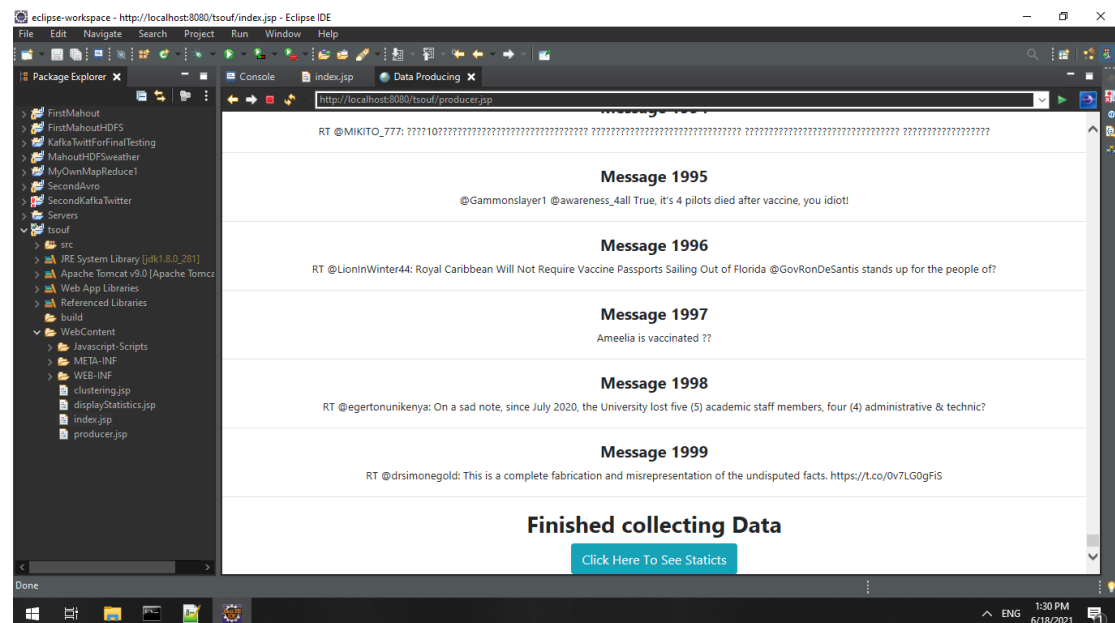
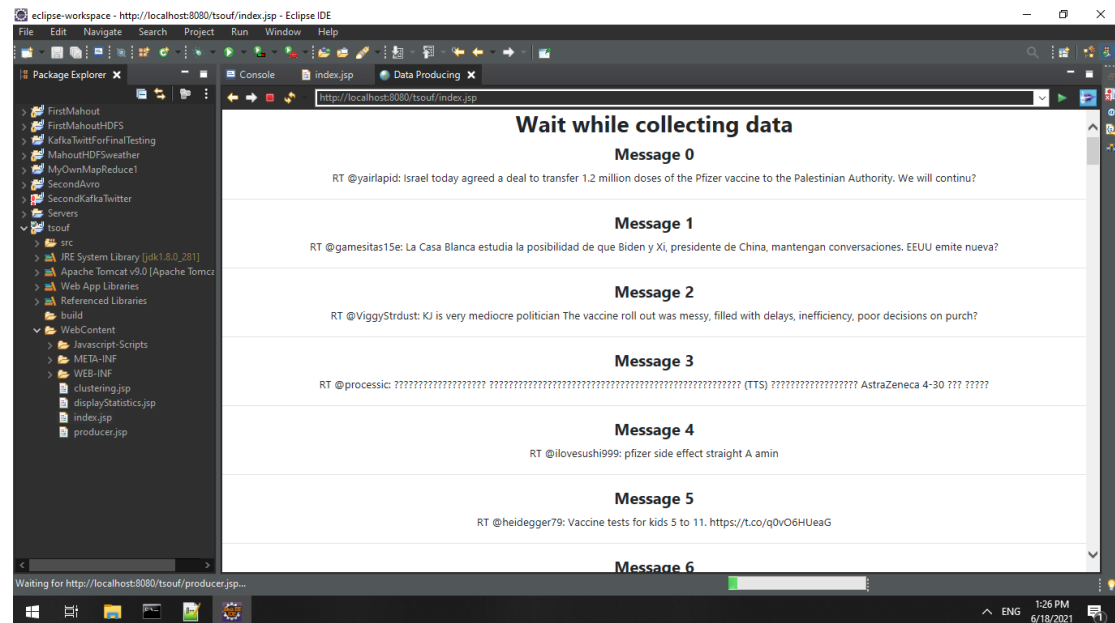
Για να τρέξει κάποιος την εφαρμογή πρέπει :

- 1) Να ανοίξει στο eclipse τον φάκελο με τα αρχεία της εφαρμογής μου (tsouf) .
- 2) Να έχει εγκαταστήσει σωστά τον Tomcat .
- 3) Να τρέχουν τα hdfs , yarn , zookeeper , kafka .
- 4) Να τρέχει το flume χρησιμοποιώντας το .conf αρχείο μου (twitter5.conf)
- 5) Να φτιάξει τοπικά στο hdfs του ένα φάκελο με το όνομα finalProject_file15
- 6) Να τρέξει την σελίδα index.jsp

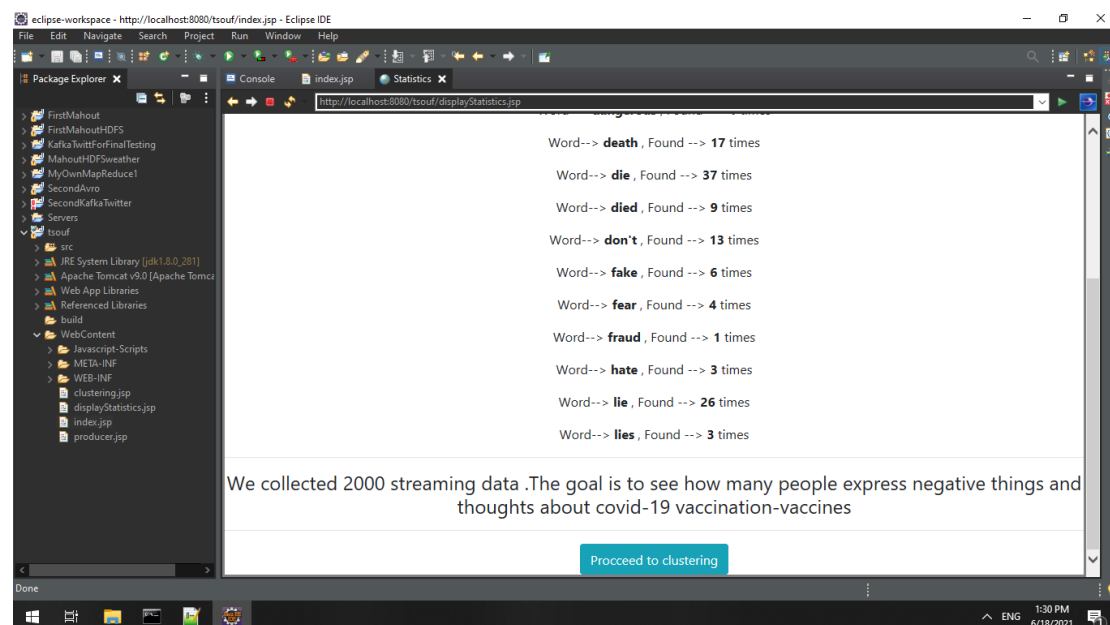
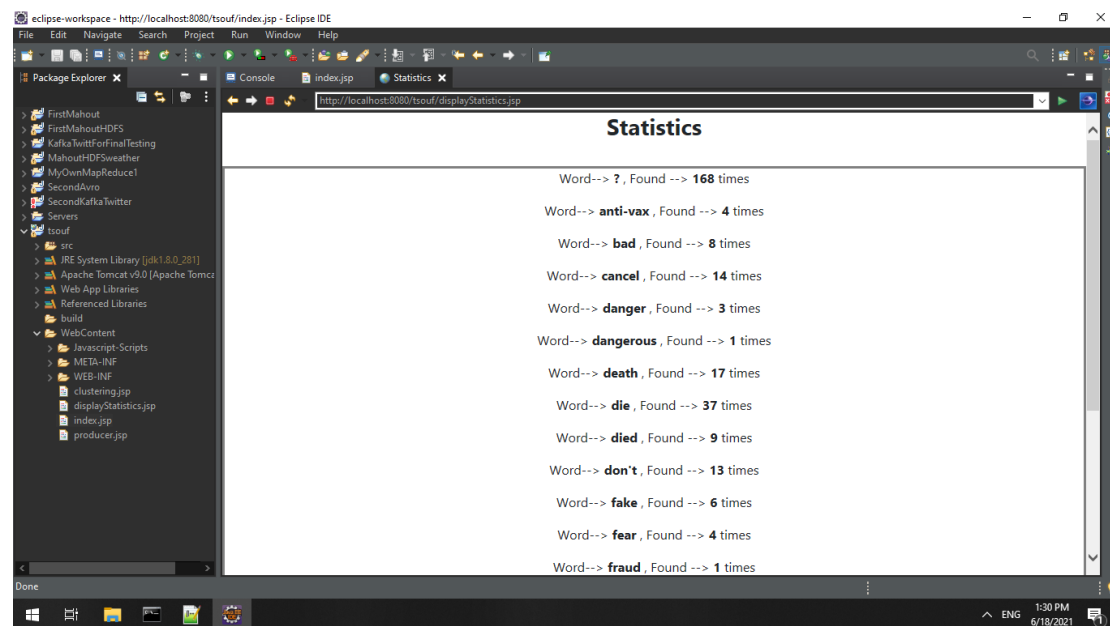
Μόλις φορτωθεί η σελίδα θα δει το παρακάτω παράθυρο .Για να ξεκινήσει η διαδικασία αρκεί να πατήσει το μπλέ κουμπί .



Στη συνέχεια περιμένει μέχρι να μαζευτούν 2000 streaming δεδομένα από το twitter με τις παραπάνω λέξεις κλειδιά.



Τώρα πατάει το μπλέ κουμπί που φαίνεται από πάνω για να δει κάποια στατιστικά .Πρόκειται για την συχνότητα εμφάνισης κάποιων συγκεκριμένων λέξεων μέσα στα συλλεχθέντα tweets .



Τώρα θα πατήσει το κουμπί που φαίνεται παραπάνω έτσι ώστε να διαχωριστούν τα μηνύματα που έχουν αρνητική χροιά από τα υπόλοιπα .

The screenshot shows the Eclipse IDE with the following components:

- Package Explorer (Left):** Displays the project structure. The 'tsouf' project is expanded, showing a 'src' folder containing 'clustering.jsp'.
- Console (Right):** Shows the URL 'http://localhost:8080/tsouf/clustering.jsp'.
- Main Editor:** Displays a large text overlay: "Everything went Ok .To see results open your hdfs UI and check folder : /KmeansOutput_finalProject_5".
- Status Bar (Bottom):** Shows the word 'Done'.

The screenshot shows the Eclipse IDE interface. On the left is the Package Explorer, displaying a project named 'FirstMahout' with a sub-package 'clustering'. The main area on the right is the Console, showing the output of the 'clustering' package. The output consists of 20 lines, each representing a cluster. A red circle highlights the text 'Note: 2 Clusters' in the console output.

ΕΝΟΤΗΤΑ 4: Αποτελέσματα Εφαρμογής

Οι λέξεις κλειδιά βάση των οποίων γίνεται η συλλογή δεδομένων , όπως αναφέρονται και στην αρχική σελίδα της εφαρμογής , είναι covid-19 , vaccine , vaccination , vaccinated , AstraZeneca , Pfizer .

Ακολουθούν παραδείγματα από τα συλλεχθέντα δεδομένα

Message 3

RT @processic: ?????????????????? ??? (TTS) ?????????????????? AstraZeneca 4-30 ??? ?????

Message 4

RT @ilovesushi999: pfizer side effect straight A amin

Message 1995

@Gammonsayer1 @awareness_4all True, it's 4 pilots died after vaccine, you idiot!

Το configuration αρχείο που δημιουργήθηκε για το εργαλείο flume είναι το ακόλουθο :

```
1 KafkaAgent.sources = Kafka
2 KafkaAgent.sinks = HDFS
3 KafkaAgent.channels = MemChannel
4
5 KafkaAgent.sources.Kafka.type = org.apache.flume.source.kafka.KafkaSource
6 KafkaAgent.sources.Kafka.kafka.bootstrap.servers = localhost:9092
7 KafkaAgent.sources.Kafka.kafka.topics = finalProject_file15
8 KafkaAgent.sources.Kafka.kafka.consumer.group.id = flume
9 KafkaAgent.sources.Kafka.interceptors = il
10 KafkaAgent.sources.Kafka.interceptors.il.type = timestamp
11 KafkaAgent.sources.Kafka.kafka.consumer.timeout.ms = 100
12 KafkaAgent.sources.Kafka.batchSize = 100
13
14 KafkaAgent.sinks.HDFS.type = hdfs
15 KafkaAgent.sinks.HDFS.hdfs.path = hdfs://localhost:9000/finalProject_file15/
16 KafkaAgent.sinks.HDFS.hdfs.fileType = DataStream
17 KafkaAgent.sinks.HDFS.hdfs.writeFormat = Text
18 KafkaAgent.sinks.HDFS.hdfs.batchSize = 1000
19 KafkaAgent.sinks.HDFS.hdfs.rollSize = 0
20 KafkaAgent.sinks.HDFS.hdfs.rollCount = 1000
21
22
23 KafkaAgent.channels.MemChannel.type = memory
24 KafkaAgent.channels.MemChannel.capacity = 10000
25 KafkaAgent.channels.MemChannel.transactionCapacity = 1000
26
27 KafkaAgent.sources.Kafka.channels = MemChannel
28 KafkaAgent.sinks.HDFS.channel = MemChannel
```

Όπως φαίνεται ,και όπως ζητάει η άσκηση ,source είναι το kafka ενώ sink είναι το hdfs .

Το MapReduce job που έγινε στα συλλεχθέντα δεδομένα κάνει το εξής ...

Ψάχνει σε κάθε tweet για κάποιες συγκεκριμένες λέξεις που υποδεικνύουν αρνητικότητα-προβληματισμό-φόβο και βρίσκει την συχνότητα εμφάνισης των λέξεων αυτών μέσα στα συλλεχθέντα δεδομένα .Το αποτέλεσμα το εμφανίζει η εφαρμογή στην διεπαφή της Statistics.jsp .

Κατά την προεπεργασία των δεδομένων , διαχωρίζονται όσα περιέχουν αρνητικότητα-προβληματισμό-φόβο από τα υπόλοιπα .

Τα συλλεχθέντα δεδομένα και τα αποτελέσματα της εφαρμογής αποδεικνύουν πως τα περισσότερα tweets που γίνονται και περιέχουν τις προαναφερθέντες λέξεις κλειδιά , έχουν ύψος αρνητικό και περιέχουν προβληματισμό ή και φόβο .